

Anatomy-Aware Measurement of Segmentation Accuracy

H.R.Tizhoosh^a and A.A. Othman^b

^aKIMIA Lab, University of Waterloo, 200 University Avenue West, Waterloo, Canada

^bDept. of Information Systems, Computers and Informatics, Suez Canal University, Egypt

ABSTRACT

Quantifying the accuracy of segmentation and manual delineation of organs, tissue types and tumors in medical images is a necessary measurement that suffers from multiple problems. One major shortcoming of all accuracy measures is that they neglect the anatomical significance or relevance of different zones within a given segment. Hence, existing accuracy metrics measure the overlap of a given segment with a ground-truth without any anatomical discrimination inside the segment. For instance, if we understand the rectal wall or urethral sphincter as anatomical zones, then current accuracy measures ignore their significance when they are applied to assess the quality of the prostate gland segments. In this paper, we propose an anatomy-aware measurement scheme for segmentation accuracy of medical images. The idea is to create a “master gold” based on a consensus shape containing not just the outline of the segment but also the outlines of the internal zones if existent or relevant. To apply this new approach to accuracy measurement, we introduce the anatomy-aware extensions of both Dice coefficient and Jaccard index and investigate their effect using 500 synthetic prostate ultrasound images with 20 different segments for each image. We show that through anatomy-sensitive calculation of segmentation accuracy, namely by considering relevant anatomical zones, not only the measurement of individual users can change but also the ranking of users’ segmentation skills may require reordering.

1. DESCRIPTION OF PURPOSE

Firefighters battling to extinguish a burning city block manage to put out the flames in 95% of the empty buildings. Many residents, however, die in the remaining 5% of the buildings.

What would we feel about the performance of those firefighters if this horrible scenario were real news? Does the number “95%” really mean anything? Wouldn’t we have preferred to let the 95% of empty buildings simply burn down, and instead, focus on those 5% with people living in them? This firefighting metaphor should illustrate the magnitude of the problem when we deal with the measurement of accuracy of organ, tumor and tissue segments in medical applications. Generally, we do focus on the whole segment without paying attention to any anatomically or pathologically significant zones inside the segment. Accuracy and its measurement is a very challenging topic in medical image analysis. Often, one can speak of accuracy when there exists a reference line, a benchmark instance, against which the current estimate or guess can be compared. We usually call this reference either “ground-truth” or, sometimes rather loosely, “gold standard” images. Ground-truth images are manual delineations created by the medical expert (e.g., radiologist, oncologists). The results of any segmentation algorithm, automated or not, can then be quantified via comparison with this ground-truth image. The accuracy of manual delineations can be measured against consensus segments among multiple experts (gold standard image). Hence, algorithms are accurate if their segments do overlap with what experts expect. That is the case in all validation procedures when we test the performance of software algorithms or the quality of manual delineations. In other words, we treat all pixels of a segment in the same way although, in many clinical cases, there are clearly different zones that are of lower or higher significance for the task at hand. As an example, when we are segmenting prostate glands for radiation treatment, the *rectal wall* is a critical zone for which the segment should exhibit highest accuracy possible. Another example is when we examine breast ultrasound lesions for diagnostic purposes. Here when the mass is mostly segmented correctly but some “spiculations” are missed, this can completely change the lesion classification based on BI-RADS guidelines.

E-mails: tizhoosh@uwaterloo.ca, a.othman@ci.suez.edu.eg

Our idea is to establish a zone-sensitive, or anatomy-aware accuracy measurement that can take into account anatomical or pathological a-priori knowledge and incorporate it into the accuracy measurement.

2. THE METHODS

There is a vast literature on evaluation of segmentation results.¹⁻⁶ The problem of validating the segmentation accuracy in medical image analysis is apparently that we look at the entire segment without any internal discrimination, meaning that some important zones inside the segment are completely ignored. What is the solution? It seems that we cannot develop any solution unless those “significant zones” inside the segment are defined prior to the calculation. But that means we have to ask the medical expert to highlight the zones in every segment, and this can be a very tedious task and hence an infeasible requirement. Keeping in mind that ground-truth segments by at least one expert must be available for any type of accuracy measurement, we cannot put additional burden of delineating the zones in individual ground-truths on the expert. So what is the solution?

The zones have to be highlighted in a “master shape”, a general or statistical shape that represents the expected shape appearance of the organ or tumor. Of course, such an approach can only address the cases with more or less regular shapes, e.g., organs and compact masses such as cysts and nodules. As well, it would need to be done only once in order to not create additional work for the clinical experts. A master shape with zones inside would then constitute a “master gold”. Every time that we have a segment and corresponding ground-truth, we can map the zones from the master gold to the current ground-truth and subsequently to the segment. This finally enables us to perform zone-sensitive accuracy measurements provided we also have some zone-sensitive accuracy measures (if we extend existing ones to become aware of zonal anatomy within the segment) to capture the compound accuracy. The outline of this idea is illustrated in Figure 1.

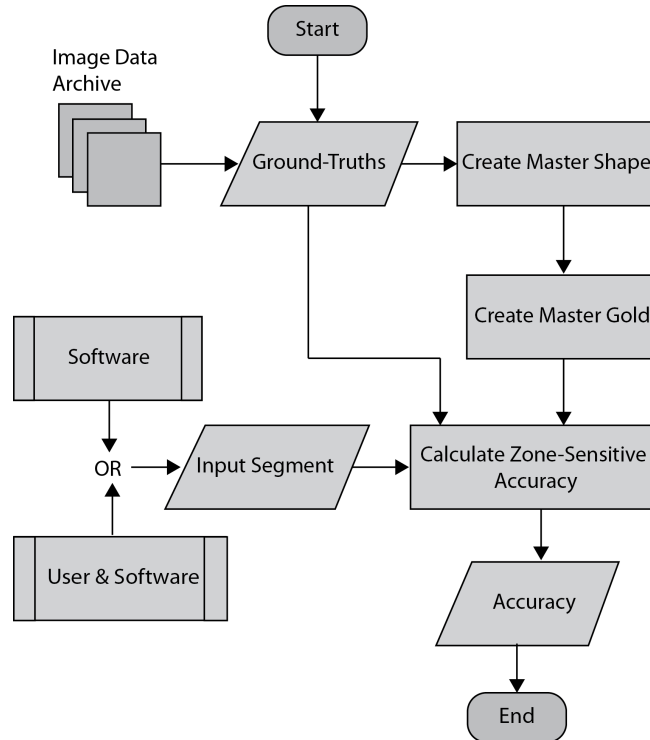


Figure 1: In addition to segments and ground-truths, a master gold should be created to calculate the zone-sensitive accuracies. The master gold depicts a generic shape, called master shape, with defined zones. The segment can come from experts and/or software. Ground-truth images come from one or multiple experts.

Measuring accuracy of segmentation is generally possible if a ground-truth segment is available. This is most of the time a manual segmentation by an expert, against which the accuracy of any segments can be measured.

If there are several manual segmentations by multiple experts available for the same image, then one may build a *consensus contour* to serve as gold standard image.

Given the segment S and the ground-truth G , the Jaccard index $J(S, G)$, sometimes called the area overlap and occasionally called Tanimoto index, can be calculated as follows:⁷

$$J(S, G) = \frac{|S \cap G|}{|S \cup G|}. \quad (1)$$

Given the segment S and the ground-truth G , the Dice coefficient $D(S, G)$ can be calculated as follows:⁸

$$D(S, G) = \frac{2|S \cap G|}{|S| + |G|}. \quad (2)$$

One can show that $J = D/(2 - D)$ and $D = 2J/(1 + J)$, hence $J < D$. It is obvious that S can come from an algorithm in which case G is the ground-truth from one or multiple users. As well, S can be manual delineation by an expert whereas G is then gold standard as consensus among multiple experts. For instance, when segmenting the prostate gland, one has to actually pay more attention to some specific zones such as the rectal wall, neurovascular bundle and urethral sphincter (Figure 2, left). In many cases, a segment may have a large overlap with the ground-truth but may not be accurate enough in significant zones (Fig. 2, right). The accuracy of such segments should be penalized according to the zonal accuracy.

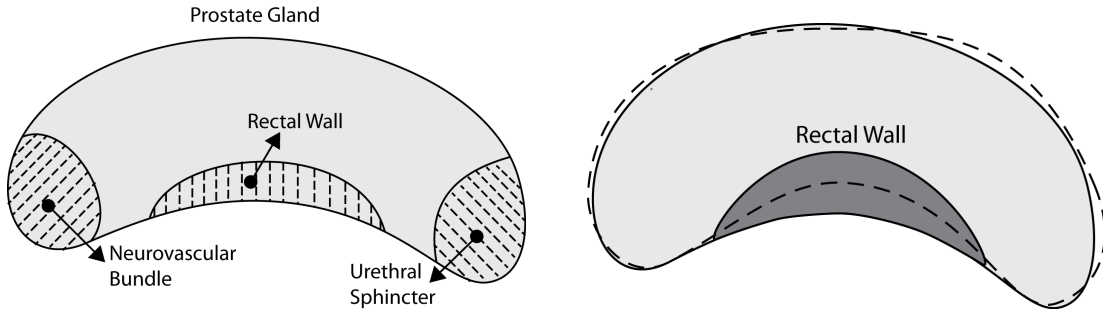


Figure 2: Significant zones within the prostate gland are generally ignored by existing accuracy measures (left). Hence, segments (dashed outline) may receive high accuracy values even though they miss a considerable portion of the rectal wall (right, dark gray).

Any anatomy-aware accuracy measure A^* with higher emphasis on zonal accuracy should hence be the extended version of an existing accuracy measure A (for instance, Jaccard or Dice) when the zonal accuracy A_Z is given and a convex combination can be built with

$$A^* = \alpha A + (1 - \alpha) \times A_Z, \quad (3)$$

where $\alpha \in [0, 1]$. Of course, if there are N_Z zones, then a representative zonal accuracy among the accuracies $A_{Z_1}, A_{Z_2}, \dots, A_{Z_{N_Z}}$ should be calculated. One may, conservatively chose

$$A^* = \alpha A + (1 - \alpha) \times \min_i A_{Z_i}. \quad (4)$$

However, the trade-off value α that determines how significant the zones are relative to the overall segment would pose another adjustment problem which is clearly not desirable. But to further investigate the establishment of a new accuracy measure, let us look at the extreme values for such trade-off parameter. In case $\alpha \rightarrow 1$ the influence of the zonal accuracy, expectedly, disappears. For $\alpha \rightarrow 0$ the zonal accuracies become dominant. However, this indicates a problem that in case the accuracy of overall segment is not high enough it would not be meaningful to pay attention to zonal accuracies. Considering these thoughts, we can establish

$$A^* = A^2 + (1 - A) \times \min_i A_{Z_i} \quad \text{if } A \geq A_{\min}, \quad (5)$$

where A_{\min} is the minimum required segment accuracy for the application at hand. For instance, if the expert/software sets $A_{\min} = 75\%$ that means the zonal accuracies will only be considered via A^* if the overall segment accuracy is at least 75%. Depending on the critical role of segmentation, any segment with $A < A_{\min}$ may be rejected as unacceptable.

Hence, to make Dice coefficient anatomy-aware, one may use

$$D_1^* = D^2 + (1 - D) \times \min_i D_{Z_i}. \quad (6)$$

Or alternatively, one may modify the core definition of the Dice coefficient to incorporate zones (TP=true positive, FP=false positive, FN=false negative):

$$D_2^* = \frac{2(\sum_i^n \text{TP} + \sum_i^{N_Z} \text{TP}_{Z_i})}{\sum_i^n (2\text{TP} + \text{FP} + \text{FN}) + \sum_i^{N_Z} (2\text{TP}_{Z_i} + \text{FP}_{Z_i} + \text{FN}_{Z_i})}. \quad (7)$$

Analogously, the Jaccard index can be extended as follows:

$$J_1^* = J^2 + (1 - J) \times \min_i J_{Z_i}. \quad (8)$$

The Jaccard extension can also occur by changing the core definition:

$$J_2^* = \frac{\sum_i^n \text{TP} + \sum_i^{N_Z} \text{TP}_{Z_i}}{\sum_i^n (\text{TP} + \text{FP} + \text{FN}) + \sum_i^{N_Z} (\text{TP}_{Z_i} + \text{FP}_{Z_i} + \text{FN}_{Z_i})}. \quad (9)$$

Extracting the Master Shape (Algorithm 1) – In order to calculate the extended accuracy measures, one apparently needs a very different approach to segmentation evaluation. Using existing ground-truth images G_i , we calculate a general (master) shape M_S . In addition to a desired minimum accuracy A_{\min} , the expert has to determine the number of zones N_Z . As well, the threshold T_{shape} needs to be set which determines the consensus level for thresholding the accumulated ground-truths (line 3, Algorithm 1) (all pixels with at least T_{shape} overlap among segments will belong to the consensus segment). One may use algorithms like STAPLE,⁹ however this failed in working with a large number of segments in our experiments such that we were forced to use our simple method to extract the master shape M_S .

Algorithm 1 Extract the General Segment Shape M_S

- 1: User sets the shape threshold T_{shape} (e.g., $T_{\text{shape}} = 50\%, 60\%, \dots$).
 - 2: Load the available gold images G_1, G_2, \dots, G_n .
 - 3: Create cumulative image: $C_G \leftarrow \sum_{i=1}^n G_i$.
 - 4: Get the master shape: $M_S \leftarrow \text{Binarize } C_G \text{ with threshold } = (n \times \frac{T_{\text{shape}}}{100})$
 - 5: Save M_S .
-

Creating the Master Gold (Algorithm 2) – In a second phase, one would need to let the expert delineate N_Z zones in the master shape M_S using N_P points (clicks) per zone to create the master gold M_G . We implemented Algorithm 2 to perform this phase, however, the zones can be delineated using any available image editor. Also one has to bear in mind that the creation of the master gold is a one-time task and generally does not need to be repeated.

As soon as a master gold M_G is available, one can start calculating the accuracy of segments using the ground-truths G provided the zones depicted in M_G can be aligned with corresponding points in the i -th ground-truth G_i and the segment S_i . Whereas the master gold M_G is one image and universally available for all images, every image I_i with the segment S_i has, as usual, its own ground-truth G_i for evaluation or training purposes.

Mapping instead of Registration (Algorithm 3) – Finding the correspondent pixels in G_i and consequently in S_i , given the zonal coordinates in M_G , seems to be a typical “registration” task. However, based on our

Algorithm 2 Create the Master Gold M_G from Master Shape M_S by acquiring the zones from user

```

1: Load the master shape  $M_S$ .
2: Set the number of (clicks) points  $N_P$ 
3: for  $i = 1 : N_Z$  do
4:   for  $j = 1 : N_P$  do
5:     Ask the user to select a point  $P_i = (x_j, y_j)$ .
6:     if  $P_i$  is close to the  $M_G$  contour then
7:       Adjust  $P_i$  to be on the contour.
8:       Save  $P_i$ 
9:     else
10:      Save  $P_i$  as a middle point
11:    end if
12:  end for
13:  Use the  $N_P$  points to create a curve  $C_i$ .
14:   $Z_i \leftarrow$  Fill in the  $i$ -th zone bounded by  $C_i$  and  $M_S$  border.
15:   $M_G \leftarrow M_S + Z_i$ 
16:  Save the coordinate of the zone  $P_i = (x_j, y_j)$ .
17: end for
18: Save  $M_G$ 

```

Algorithm 3 Map Zones to the Segment (see Algorithm 5 and Figure 3)

```

1: Load the current segment  $S$  and the Master Gold  $M_G$ .
2: for  $i = 1 : N_Z$  do
3:   if the zone at the right or the left then
4:     Apply the  $x$ -values at the  $x$ -axis on  $P_N$  to calculate the  $y$ -values.
5:   else
6:     Apply the  $y$ -values at the  $y$ -axis on  $P_N$  to calculate the  $x$ -values.
7:   end if
8:   Draw a curve using  $x$  and  $y$  values.
9:   Fill in the area under the curve that belong to the segment to create the zone.
10:  Remove any part of the curve that fall out of the segment.
11:  Save the coordinates of the zone  $(x_{S_i}, y_{S_i})$ .
12: end for

```

experimental results we decided to not use registration algorithms for this purpose. The non-rigid registrations we tested were both time-consuming (which may not be a critical drawback) and inaccurate. Whereas one may use a specific registration algorithm in context of a familiar segmentation task, we do provide a quasi-non-rigid mapping procedure that is very fast, due to its simplicity, and can handle small irregularities quite easily. For this, first we do fix some points on the contour of the master gold (see Algorithm 4 in Appendix) and then map them to the ground-truth (see Algorithm 5 in Appendix) and segment (Algorithm 3; see Figure 3).

3. RESULTS

Only organs and regular-shaped anomalies (cysts, nodules etc.) are considered. We further assume that there is at least one expert who has created ground-truth segments for each image and there is at least one expert who can mark anatomically meaningful zones with higher significance for segmentation. And finally we assume that the zones always touch the boundary of the segment.

3.1 Image Data: Synthetic TRUS Images

It is a challenge to validate any approach to segmentation. One has to measure the accuracy of the segment S against ground-truth images. Ideally, if we have many users available to segment images, we can build “consensus segments”, or *gold standard*, to make more reliable measurements. Of course, this is usually not feasible with

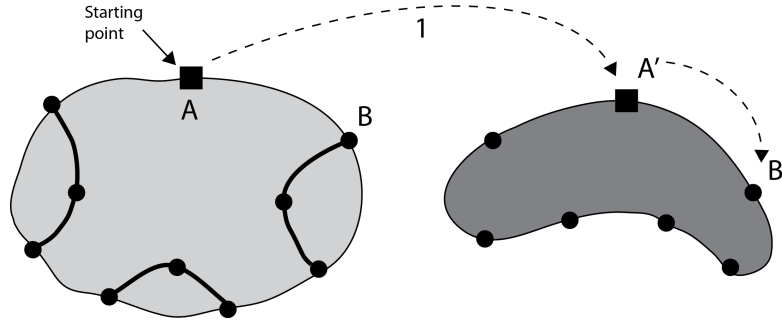


Figure 3: Point Mapping: The salient points of the zones defined in the master gold M_G (left) are mapped into the current ground-truth and segment (right).

real images, for which there is no gold standard. Hence, we generated synthetic images whose gold segments were known a priori. For this reason, we used synthetic images that simulate transrectal ultrasound (TRUS).

TRUS images of prostates may be used to both diagnose and treat prostate diseases such as cancer. Starting with a set of prostate shapes P_1, P_2, \dots, P_m , we created random segments G_i through combinations of those priors, adding noise along with random translations and rotations, and we distorted the results with speckle noise and shadow patterns. Each image I_i is thus created from its gold G_i . Consequently, we can simulate k user delineations $S_i^1, S_i^2, \dots, S_i^k$ by manipulating G_i via scaling, rotation, and morphological changes, and we can simulate edits by running active contours with variable user-simulating parameters. The variability of user delineations was simulated according to several factors: error probability ($[0, 0.05]$), anatomical difficulty ($= 0.2$ out of $[0, 1]$), and the scaling factor for morphology (from 1×1 to 21×21). The user was modelled according to the level of experience (a random number from $(0, 1]$), the user's attention (a random number from $[0, 1]$), and the user's tendencies in terms of the segment size (a random number from $[-1, 1]$), whether tending to draw contours that are relatively small ($\rightarrow -1$) or large ($\rightarrow +1$).

We generated 500 images from their corresponding gold-standard images*. Furthermore, we generated 20 different segments for each image, assuming that there were 20 users. Figure 4 shows five examples of real and synthetic TRUS images. One should bear in mind that the purpose here was not to simulate the images realistically, but rather to have a base from which to generate variable segments from a perfect segment. Figure 5 shows an example of the gold segments and simulated user contours. The variability, coupled with the gold segment, is what is needed in our experiments.

3.2 Experiments

We conducted several experiments to examine the effect of employing the new accuracy measures. In the first experiments we measured the accuracy of all 10,000 segments (500 images each segmented by 20 simulated users). The accuracy measurement encompassed the conventional Jaccard index \bar{J} , the Jaccard values for the three zones \bar{J}_{Z_1} , \bar{J}_{Z_2} and \bar{J}_{Z_3} , as well as the two variations of total Jaccard accuracies for the entire segments \bar{J}_1^* and \bar{J}_2^* . These results are reported in Table 1. It is apparent the extended Jaccard values are lower than the conventional ones: $\bar{J} > \bar{J}_1^* > \bar{J}_2^*$. The selection of the best segment may change depending on the measure whereas zonal accuracies show a more pronounced shift. In particular, if one chooses \bar{J}_{Z_2} (zone 2) as a base, the results may have a different impact with respect to the quality of the segments. Similar results were observed for Dice coefficient D and its anatomy-aware version D^* .

As a subset of the experiments, we randomly selected 50 images and 10 simulated users to examine some details (see Table 2). Both versions of anatomy-aware Jaccard deliver lower accuracies for any given user. Whereas J_1^* is on average 10% lower, J_2^* is about 16% lower. The zone 1 seems to be the most difficult zone for almost all users. However, some users (e.g., users 1, 3, 4 and 6) appear to be more challenged with the zone 3. Users 8 and 9 are the best users ($\bar{J} = 87$ and 86, respectively). Their performance, however, is quite low when

*All images and their segments are available online: <http://tizhoosh.uwaterloo.ca/>

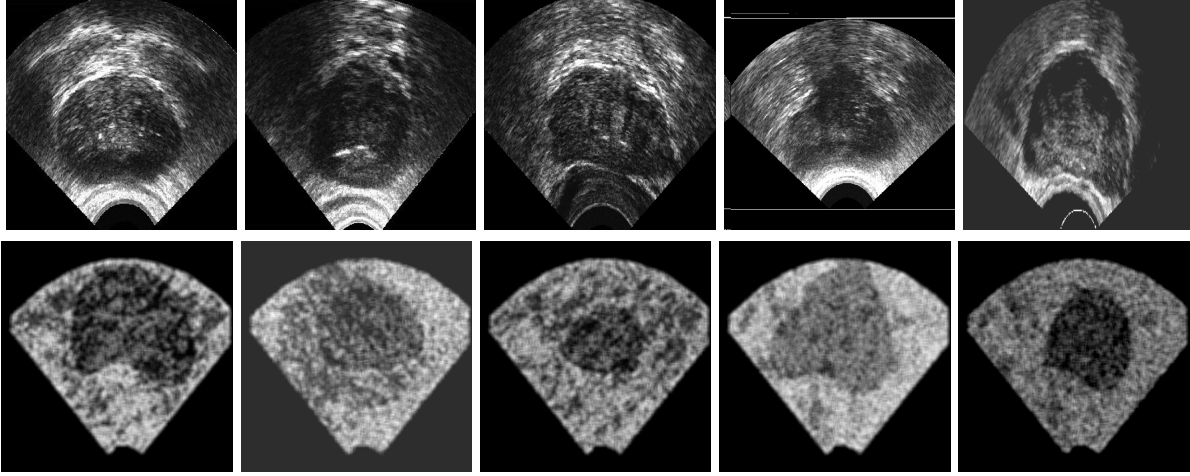


Figure 4: Sample TRUS images (top) and simulated images (bottom).

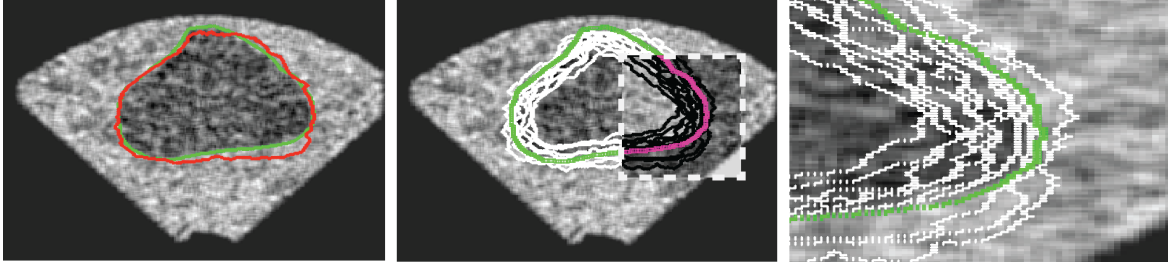


Figure 5: Left: Sample image with gold segment and consensus contour; Middle: Simulated user segments with the gold contour; Right: The inverted region (middle) magnified to show the variability.

Table 1: All results for images with $J > 75\%$. Highest accuracies are highlighted for each measure.

User	\bar{J}	\bar{J}_{Z_1}	\bar{J}_{Z_2}	\bar{J}_{Z_3}	\bar{J}_1^*	\bar{J}_2^*
1	87± 6	64± 20	83± 12	61± 22	82± 10	74± 11
2	82± 5	48± 17	78± 13	50± 19	74± 9	67 ± 7
3	83± 4	63± 13	75± 11	53± 19	77± 7	67± 8
4	85± 5	64± 13	76 ± 11	57± 20	79± 7	70± 9
5	80± 4	45± 13	75± 10	43± 15	71± 6	65± 6
6	81± 4	59± 14	72± 10	49± 16	74± 6	64± 7
7	80± 4	45± 14	74± 12	44± 16	71± 7	65± 5
8	86± 7	54± 23	82± 14	64± 22	80± 11	73± 10
9	88± 5	66± 17	81± 11	63± 21	83± 8	74± 9
10	81± 4	59± 13	72± 11	49± 18	74± 7	64± 8
11	87± 6	64± 20	83± 12	61± 22	82± 10	74± 11
12	82± 5	48± 17	78± 13	50± 19	74± 9	67± 7
13	83± 4	62± 13	75± 11	53± 19	77± 7	67± 8
14	85± 5	64± 13	76± 11	57± 20	79± 7	70± 9
15	80± 4	45± 13	75± 10	43± 15	71± 6	65± 5
16	81± 4	59± 14	72± 10	49± 16	74± 6	64 ± 7
17	80± 4	45± 15	75± 11	44± 16	71± 7	65± 5
18	86± 7	54± 23	82± 14	64± 22	80± 11	73 ± 10
19	88± 5	66± 17	81± 11	63± 21	83± 8	74± 9
20	81± 4	59± 13	72± 11	49± 18	74± 7	64 ± 8

segmenting the zone 1 ($\bar{J}_{Z_1} = 59$ and 66, respectively). Their performance seems to be more plausibly captured by the first anatomy-aware measure ($\bar{J}_1^* = 81$ and 82, respectively) which also favors user 9 instead of user 8. The second anatomy-aware measure appears to be very conservative ($\bar{J}_2^* = 74$ and 72, respectively). Both standard deviation and variance illustrate that user variability is amplified by variability in zones 1 and 3. The first anatomy-aware measure, \bar{J}_1^* , seems to more pronouncedly quantify the user variability. Table 3 shows how the ranking of users change when we base our evaluations upon anatomy-aware measures. Apparently, the ranking of users with excellent segmentation skills may not change much. In contrast, considerable shift in ranking can be observed when the user skills is rather average. For users with high Jaccard value, the ranking does not seem to change (users 3, 4, 8 and 9). Users with poor segmentation skills (user 10) does not seem either to change their ranking. For users with “average” skills (Jaccard values around 60%-70%), the ranking may considerably change if we use anatomy-aware Jaccard (gray rows in Table 3).

Table 2: Accuracy measurements via conventional Jaccard (first column), the defined three zones (gray columns), and the two anatomy-aware versions of Jaccard (last two columns).

User	\bar{J}	\bar{J}_{Z_1}	\bar{J}_{Z_2}	\bar{J}_{Z_3}	\bar{J}_1^*	\bar{J}_2^*
1	69	46	57	44	60	51
2	74	35	68	48	61	57
3	79	59	66	56	72	62
4	78	57	65	56	71	61
5	70	30	59	42	57	53
6	72	50	59	48	64	54
7	72	37	66	39	59	55
8	87	59	80	74	81	74
9	86	66	76	71	82	72
10	57	12	45	18	36	39
STDV	9	17	10	16	13	10
variance	69	248	88	234	161	94

Table 3: Ranking of segmentation skills of simulated users based on different accuracy measures.

Rank	\bar{J}	\bar{J}_{Z_1}	\bar{J}_{Z_2}	\bar{J}_{Z_3}	\bar{J}_1^*	\bar{J}_2^*
1	8	9	8	8	9	8
2	9	8	9	9	8	9
3	3	3	2	3	3	3
4	4	4	3	4	4	4
5	2	6	7	2	6	2
6	6	1	4	6	2	7
7	7	7	6	1	1	6
8	5	2	5	5	7	5
9	1	5	1	7	5	1
10	10	10	10	10	10	10

4. CONCLUSIONS

We introduced the novel idea of anatomy-aware accuracy measures. Extending commonly used measures such Jaccard index and Dice coefficient to anatomy-sensitive schemes is proposed by designing multiple necessary algorithms. Among others, the concept of “master gold” is introduced which is necessary for implementation of any anatomy-aware accuracy measurement. Anatomy-sensitive accuracy measurement appears to provide more insight into the challenges of medical image segmentation. By considering anatomical zones within segments, we

may be able to develop a better understanding of contouring skills of users. As well, anatomy-aware accuracy measures seem to provide a more realistic qualification of inter-observe variability. And finally, anatomy-aware measures can be used to improve the performance of trainable segmentation accuracy.^{10–12}

Contributions

The extensions of Jaccard and Dice measures to their zonal versions have been designed by the first author. As well, the image simulation to generate test data was designed and implemented by the first author. The second author has conducted all experiments and generated all results. The paper has been written by the first author.

Acknowledgements

The authors would like to thank Dr. Masoom Haider (Sunnybrook Research Institute, Toronto) for some early discussions and advice with respect to the anatomy of the prostate gland. Also Dr. Haider provided us with some insight into the nature of the accuracy problem for the prostate gland. As well, the authors would like to thank Dr. Farzad Khalvati (Dept. of Medical Imaging, University of Toronto) for some initial elaborations on how experiments should be conducted.

This project was funded by the Natural Sciences and Engineering Research Council of Canada (NSERC) in form of a Discovery Grant.

REFERENCES

- [1] Popovic, A., de la Fuente, M., Engelhardt, M., and Radermacher, K., “Statistical validation metric for accuracy assessment in medical image segmentation,” *International Journal of Computer Assisted Radiology and Surgery* **2**(3-4), 169–181 (2007).
- [2] Shepherd, T., Teras, M., Beichel, R. R., Boellaard, R., Bruynooghe, M., Dicken, V., Gooding, M. J., Julyan, P. J., Lee, J. A., Lefevre, S., Mix, M., Naranjo, V., Wu, X., Zaidi, H., Zeng, Z., and Minn, H., “Comparative study with new accuracy metrics for target volume contouring in pet image guided radiation therapy,” *IEEE Transactions on Medical Imaging* **31**(11), 2006–2024 (2012).
- [3] Zou, K. H., Warfield, S. K., Bharatha, A., Tempany, C. M. C., Kaus, M. R., Haker, S. J., Wells III, W. M., Jolesz, F. A., and Kikinis, R., “Statistical validation of image segmentation quality based on a spatial overlap index,” *Academic Radiology* **11**(2), 178–189 (2004).
- [4] Correia, P. and Pereira, F., “Objective evaluation of relative segmentation quality,” in [*IEEE International Conference on Image Processing*], **1**, 308–311 (2000).
- [5] Chang, H. ., Zhuang, A. H., Valentino, D. J., and Chu, W. ., “Performance measure characterization for evaluating neuroimage segmentation algorithms,” *NeuroImage* **47**(1), 122–135 (2009).
- [6] Udupa, J. K., LeBlanc, V. R., Zhuge, Y., Imielinska, C., Schmidt, H., Currie, L. M., Hirsch, B. E., and Woodburn, J., “A framework for evaluating image segmentation algorithms,” *Computerized Medical Imaging and Graphics* **30**(2), 75–87 (2006).
- [7] Jaccard, P., “The distribution of the flora in the alpine zone,” *The New Phytologist* **XI**(2), 37–50 (1912).
- [8] Dice, L., “Measures of the amount of ecologic association between species,” *Ecology* **26**(3), 299–302 (1945).
- [9] Warfield, S. K., Zou, K. H., and Wells, W. M., “Simultaneous truth and performance level estimation (staple): An algorithm for the validation of image segmentation,” *IEEE Transactions on Medical Imaging* **23**(7), 903–921 (2004).
- [10] Othman, A., Tizhoosh, H., and Khalvati, F., “EFIS – evolving fuzzy image segmentation,” *Fuzzy Systems, IEEE Transactions on* **22**(1), 72–82 (2014).
- [11] Othman, A. and Tizhoosh, H., “N-cuts parameter adjustment using evolving fuzzy inferencing,” in [*Fuzzy Systems (FUZZ), 2013 IEEE International Conference on*], 1–6 (2013).
- [12] Sahba, F., Tizhoosh, H., and Salama, M., “Application of opposition-based reinforcement learning in image segmentation,” in [*Computational Intelligence in Image and Signal Processing, 2007. CIISP 2007. IEEE Symposium on*], 246–251 (2007).

Appendix

Algorithm 4 Determine salient contour points on M_G outline for mapping

```

1: Load the master gold  $M_G$ .
2: Load the points  $P_1, P_2, \dots, P_{N_Z}$ .
3: Copy the border points  $P$  into  $P_B$ .
4: Get the number of border points  $N_{P_B}$ .
5: % — Calculate relative border distances —
6: Determine a starting point  $C$  on  $M_G$ 's contour (see Figure 3).
7: Get  $M_G$ 's contour,  $(X_C, Y_C)$ , starting from  $C$ .
8: Get the segment length  $L = \max(X_C) - \min(X_C)$ 
9: Get the segment width  $W = \max(Y_C) - \min(Y_C)$ 
10: for  $i = 1 : N_{P_B}$  do
11:   Calculate the distance  $D(i, 1)$  from  $C$  to  $P_B(x_i, y_i)$ .
12:   Normalize the distance  $D(i, 2) = D(i, 1)/|X_C|$ .
13: end for
14: % — Calculate relative internal distances —
15: Copy the middle points from  $P$  to  $P_M$ .
16: for  $i = 1 : N_Z$  do
17:   Determine a starting point  $C_Z$  on the border of the  $i$ -th zone between the zone end points.
18:   Calculate the distance  $D_W(i, 1) = \|C_{Z_i}, P_{M_i}\|$ .
19:   % — Normalize the distance —
20:   if the zone on the right or on the left then
21:      $S = W$ 
22:   else
23:      $S = L$ 
24:   end if
25:    $DW(i, 2) = DW(i, 1)/S$ .
26: end for
27: Add  $C_Z$  to  $P$ .
28: Save  $D, D_W, P$ .

```

Algorithm 5 Map zones to the ground-truth G

```

1: Load  $P, D, D_W$ 
2: Read the current ground-truth image  $G$ .
3: Determine a starting point  $C_G$  on  $G$ 's contour (see Figure 3).
4: Get  $G$ 's contour,  $(X_G, Y_G)$ , starting from  $C_G$ .
5: get the length  $L_G$  and the width  $W_G$  of  $G$ .
6: Calculate the distance  $D_G$  from  $C_{G_Z}$  to the suggested zone border points on  $G$ :  $D_G = D(:, 2) \times \text{length}(X_G)$ .

7: Calculate the point at the border of each zone  $P_{GB}$  on  $G$  using  $D_G$  and  $(X_G, Y_G)$ :  $P_{GB} = [X_G(D_G)Y_G(D_G)]$ .

8: Calculate the centre points at the border of each zone  $C_{G_Z}$  using  $P_{GB}$ .
9: Calculate  $S_G$  the same way as  $S$ .
10: Calculate the distance  $D_{G_W}$  from  $C_{G_Z}$  to the middle point of the zone:  $D_{G_W} = D_W(:, 2) \times S_G$ 
11: Calculate the point at the curve of each zone  $P_{G_M}$  using  $C_{G_Z}$  and  $D_{G_W}(:, 2)$ .
12: Using  $P_{GB}$  and  $P_{G_M}$ , draw the curve of the zone.
13: Save the coordinates of the zone  $(x_{G_i}, y_{G_i})$ .
14: Save the polynomial parameters  $P_N$  used for drawing the curve.

```
